# Data Mining Tools and Technique

**Yogesh Shah**
Student, Department of Masters of Computer Applications
Late Bhausaheb Hiray S. S. Trust's Hiray Institute of Computer Application, Mumbai, India

**Abstract**: *Data mining is one of the research areas that is most suited to computer applications among the various types of data analysis. Mining the web will be the main topic of this essay. This review study shows a comprehensive examination of the various web mining methods and instruments. The set of tools needed for fulfilling the promise is mining the web. It is the use of data mining tools to derive knowledge from the usage, structure, and content of websites. The definition of web mining is discussed in this essay. A wide range of data types, including scientific, environmental, financial, and mathematical data, are everywhere in the world. Due to the rapid growth of data in the current era of the internet and information exchange, manual analysis, classification, and summarization are no longer feasible. This research discusses the basic principles of data mining as well as contemporary research on the topic in an effort to create novel methods for incorporating handling uncertainty into the use of data mining*

**Keywords:** Data, Mining, clustering, classification

## I. INTRODUCTION

Most internal auditors, especially those that work in customer-oriented groups, have heard of data mining and what it can do for an organisation, including reducing the cost of getting prospective customers and increasing the sales of new products and services. However, whether you are a new internal auditor or a seasoned veteran looking for a refresher, possessing an excellent grasp of what data mining does and the various data mining tools and techniques available for utilisation will enhance the auditing process and business operations all around.
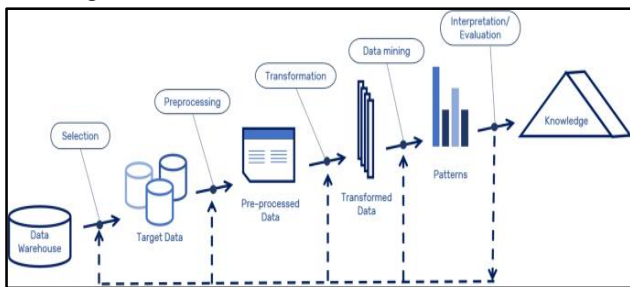
## II. WHAT IS DATA MINING?

In the most basic sense, data mining automates the recognition of using specific methods and algorithms to examine recent and historical data, finding patterns in a database that can then examined in order to forecast future trends. Data mining technologies enable organisations to make proactive, knowledge-driven decisions and provide answers to issues that were previously too time-consuming to address since they forecast future trends and behaviours by scanning databases for hidden patterns. Data mining is not particularly new; for many years, statisticians examined data and offered business estimates using comparable manual techniques. Adjustments to data mining methods, but have made it possible for businesses to gather, examine, and use data in new ways. The initial modification concerned the acquisition of fundamental data. Prior to moving away from ledgers and other paper-based records in favour of computer-based systems, managers had to wait for personnel to put the puzzle together in order to understand how the business was performing or how recent performance periods compared to earlier performance periods. Companies were able to begin as soon as they began gathering and storing basic data in computers.

Data mining methods have also been impacted by changes in data access, where there has been more empowerment and integration, especially in the last 30 years. The development of middleware, protocols, and other approaches that allow data to be transported effortlessly among programmes and other machines, as well as the arrival of microcomputers and networks, allowed businesses to link specific data questions together. For example, the development of data warehousing and decision support systems has allowed businesses to expand queries from "What was the total number of sales in New South Wales last April?" to "What was the total number of sales in New South Wales this April?" is likely to happen to sales in Sydney next month, and why?"However, the primary distinction between current data mining initiatives and earlier ones is that organisations now have access to more information. The enormous volumes of Companies frequently use data mining software to swiftly evaluate

huge amounts of data and look for trends in the information they collect. Users can control the data analysis's outcome using the parameters they choose, adding value to company strategies and objectives. The data mining programme will produce all permutations or combinations, regardless of their relevance, if these parameters are not provided.

This final aspect is important for internal auditors to understand: because data mining programmes lack the human intuition to distinguish between a useful data correlation and an irrelevant connection, users must review the outcomes of mining operations to make sure they offer the information required. For instance, it can be important to know that people who fail on loans frequently provide a bogus address, but it might not be necessary to know that they have blue eyes. Therefore, auditors should keep an eye on whether judgements based on data mining exercises are sensible and rational, especially when those decisions are used as input for other systems or processes. The many security facets of data mining programmes and processes must also be taken into account by auditors. An outsider who infiltrates the competitor organization's computer system and utilises a data mining tool on the obtained information may be able to leverage the critical customer information revealed by a data mining operation to their advantage.



## III. DATA MINING TOOLS

Companies interested in using data mining tools have two choices: either they may purchase mining programmes intended for current hardware and software platforms and integrate them into just-released products and systems when they go online, or they can develop their own distinctive mining solution. The output of a data mining operation, for instance, is frequently fed into another computer system, such a neural network, to potentially boost the value of the mined data. This is because the data mining tool collects the data, whilst the second programme (such the neural network) makes decisions depending on the data obtained.

There are many kinds of data mining tools accessible, each with advantages and disadvantages of its own.

Internal auditors must be knowledgeable about the various categories of data mining technologies that are available and urge the acquisition of a tool that meets the organization's present investigative requirements. As early as feasible in the project's lifecycle—possibly even in the feasibility study—this should be considered in account.

Traditional data mining tools, dashboards, and text mining tools are the three categories into which the majority of data mining tools fall. Here is an explanation of each.

### 3.1 Traditional Data Mining Tool

By utilising a variety of intricate algorithms and methodologies, traditional data mining programmes assist businesses in identifying data patterns and trends. Some of these tools are installed on the desktop to track data and identify trends, while others gather data that isn't stored in a database. The majority are offered in Windows and UNIX versions, while some are only compatible with that one platform. Additionally, while some may focus on one type of database, the majority will be able to handle any database.

### 3.2 Dashboard

Dashboards are systems that are installed in computers to monitor data in databases. They display data updates and changes onscreen, frequently in the form of a chart or table, allowing the user to see how **the** company is doing. In order to see where things have changed (such as an increase in sales from the same time last year), historical data can also be referred to. Because of this feature, dashboards are simple to use and particularly appealing to managers who want to get a general sense of how the business is performing.

### 3.3 Text Mining Tools

Because the third category of data mining tools can extract data from many types of text, including plain text files and Microsoft Word and Acrobat PDF documents, it is frequently referred to as a text-mining tool. These programmes scan information and transform the chosen data into a format that is compatible with the tool's database, giving users a quick and simple way to retrieve the data without having to launch other programmes. Emails, Internet sites, music, and video files are examples of unstructured scanned content. On the other hand, structured content is data whose shape and purpose are recognised, such as information found in databases.

## IV. DATA MINING TECHNIQUE

Internal auditors can select from a variety of data analysis strategies in addition to employing a specific data technology for mining. The nearest-neighbour method, decision trees, and neural networks that are artificial are among the most often utilised techniques.
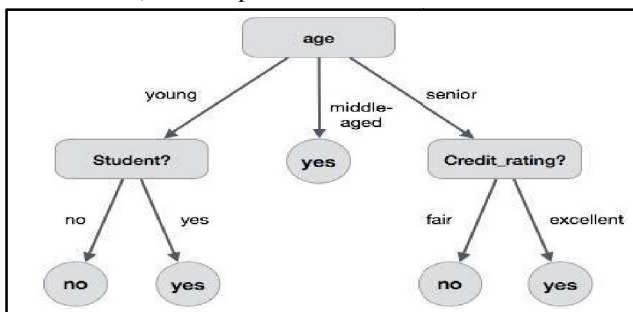
These methods each analyse data in a different way:

### 4.1 Artificial Neural Network

Artificial neural networks are non-linear, training-based predictive algorithms. Despite being effective predictive modelling techniques, part of their effectiveness comes at the expense of their usability and deployment. When evaluating records to find fraud and behaviours that may be considered fraud, for example, auditors can use them with ease. They are more effective because of their complexity when applied in circumstances where they can be repeated, like monitoring credit card transactions each month to look for irregularities.

### 4.2 Decision Tree

Decision trees are tree-shaped structures that represent decision sets. These decisions generate rules, which then are used to classify data. Decision trees are the favoured technique for building understandable models. Auditors can use them to assess, for example, whether the organization is using an appropriate cost-effective marketing strategy that is based on the assigned value of the customer, such as profit.



### 4.3 The Nearest-Neighbour Method

According to similar information in a historical dataset, the nearest-neighbour technique categorises dataset records. Using this method, auditors can specify a document that piques their attention and ask the system to look for related documents.
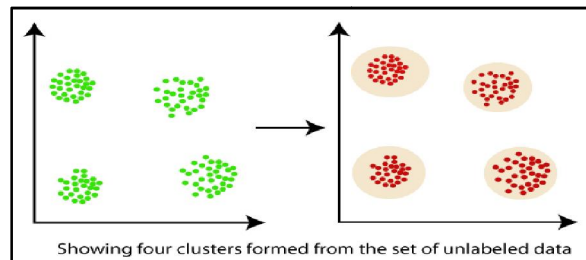
### 4.4 Association

Association (or relation) is perhaps the more well-known, well-known, and simple data mining technique. In order to find patterns, you construct a straightforward link between two or more things, frequently of the same sort. If you keep track of people's purchasing patterns, for instance, you might notice that a client always buys cream when they buy strawberries, and you might then advise them to consider purchasing cream the next time they buy strawberries.

### 4.5 Clustering

You can combine different bits of data to generate a structure opinion by examining one or more attributes or classes. In its most basic form, clustering is the process of recognising a group of related outcomes based on one or more attributes. Because it connects with other samples so you can see where the similarities and ranges agree, clustering is important for identifying different information. Clustering has two possible outcomes. You can make the assumption that a cluster exists at a specific location and then check to see whether you're right by using our identification criteria. The graph in Figure 3 provides a nice illustration. In this illustration, a sample of sales data compares the customer's age to the amount of the sale. It is logical to anticipate that those in their twenties—before marriage and having children—and their fifties and sixties—after their kids have fled the nest—will have greater disposable income.


Showing four clusters formed from the set of unlabeled data

### 4.6 Classification

By specifying several criteria to pinpoint a specific class, classification can help you develop an understanding of the different types of customers, products, or objects. For instance, by identifying several characteristics (such as the number of seats, car form, and driven wheels), you may categorise cars into various classes (such as sedan, 4x4, and convertible). You can define a brand-new car by comparing its characteristics to a definition we are familiar with. Customers can be categorised by age and social group, for example, in accordance with the same criteria. Classification can also be the output of other approaches or a feeder to them. Decision trees can be used, for instance,

to decide on a classification. Using shared qualities from many classes, clustering enables you to locate clusters.

### 4.7 Sequential Patterns

Sequential patterns, which are often used across longer-term data, are an effective way to spot trends or recurrent occurrences of comparable events. For instance, you may learn from customer data that clients buy a specific group of products together at various points during the year. Based on their frequency and previous purchases, you can use this data in a shopping basket application to automatically propose that specific items be put to a basket.

### 4.8 Predication

Prediction is a broad topic that includes everything from foreseeing the breakdown of parts or equipment to detecting fraud and even forecasting business earnings. Prediction makes use of trend analysis, classification, pattern matching, and relation when combined with other data mining approaches. You can forecast an event by looking at previous occurrences or happenings. To determine whether a transaction is fraudulent, you might use the credit card authorisation as an example and combine decision tree analysis of specific previous transactions with classification and historical pattern matches. Finding a link between booking flights to the US and transactions there suggests that the transaction is probably legitimate.

### V. EFFECTIVE WAY TO OPTIMIZE DATA MINING

Multi-database mining is becoming more crucial for efficient and informed decision-making as businesses expand to multiple locations. You may make the most of your data mining efforts by utilising the following mining strategies.

Step 1: Accurate classification is hampered by incomplete data, which also makes data mining less efficient.

Step 2: The creation, upkeep, and performance optimisation of a parallel data mining application currently demand a lot of knowledge and work.

Step 3: Combining a variety of architectural options for connecting mining with database systems is a good data mining strategy.

Step 4: It is essential to create a system that enables interactive multiple-level knowledge mining in massive relational databases and data warehouses. Online analytical processing (OLAP) must be tightly integrated

with a wide range of data mining tasks, including characterization, association, classification, prediction, and clustering.

Step 5: In the year 2011, heterogeneous systems of databases are crucial to the information industry. To stay up with the trend, data warehouses must offer data extraction from several databases.

For instance, modelling the actions of telecom organisations necessitates the use of three heterogeneous data mining programmes. First, the neural network method is used to determine the client's attribute weight from the original data. The decision tree method is then used to identify exceptional client characteristics based on attribute weight. Finally, an adaptive differentiating model based on clustering is created. Three algorithms combined together make it easier to identify excellent clients.

### VI. DATA MINING APPLICATIONS

Simply put, data mining is the computer process of examining data from many angles, dimensions, and views before classifying or analysing it to provide useful information. Any form of information, such as that found in data warehouses, relational databases, multimedia databases, spatial databases, time-series databases, and the internet, can be processed using data mining techniques.

In the knowledge-based economy, data mining provides competitive advantages. This is done by giving users the most information possible to quickly come to informed business choices despite the vast amount of data at their disposal.

Data mining has resulted in several tangible benefits in a variety of application fields. So let's talk about several data-mining applications:

**Scientific Analysis**

Every day, large amounts of data are produced through scientific models. This comprises information gathered from nuclear research facilities, information on human psychology, etc. The study of these data is possible with the aid of data mining tools. Now, we are able to collect and store more new data than we can process through analysis. An illustration of scientific analysis:

- Analysing sequences in bioinformatics
- Putting astronomical objects into categories support for medical decisions.

**Intrusion Detection System**

Any unapproved operation on a digital network can be referred to as a network intrusion. Theft of priceless

network resources is a common aspect of network invasions. The process of data mining is essential for looking for abnormalities, network attacks, and intrusions. These methods assist in choosing and enhancing pertinent and usable facts from enormous data collections. The classification of pertinent data for intrusion detection systems is aided by data mining techniques. Network traffic is alerted by the intrusion detection system to external intrusions in the system. For instance:

- Identify security breaches
- Detecting Misuse
- Anomalous Finding

### Business Transaction

Every industry of business is stored in memory forever. These business-to-business or business-to-consumer interactions are typically time-related. The effective and timely use of the data in an adequate amount of time for competitive decision-making is the most important challenge for businesses who struggle to prosper in a highly competitive market. Data mining is helpful in evaluating these commercial links, creating marketing tactics, and making decisions. Example:

- Local mailing lists
- Trading markets
- Dividing up customers
- Prediction of churn

### Health Care & Insurance

To better target high-value doctors and determine which marketing initiatives will be most effective in the months to come, pharmaceutical companies might monitor their new deals force activity and their results. In contrast, data mining in the insurance sector can help anticipate which consumers will purchase new policies, detect risky client behaviour patterns, and spot fraudulent behaviour among consumers.

- Analysis of the claims, or which medical procedures are claimed jointly.
- Find effective medical treatments for various diseases.
- predicts office visits via patient behaviour characteristics.

### Financial/Banking Sector

A credit card provider can use the huge database of consumer transaction data in its hands to pinpoint those customers most likely to be drawn to a novel credit product.

- Identifying credit card fraud.
- Recognise 'Loyal' clients.
- Obtaining customer-related information.
- Analyse consumer expenditure on credit cards.

### Transportation

Data mining can be used by a diverse transportation firm with a sizable direct sales team to pinpoint the ideal customers for its products. Information mining can be used by a sizable consumer goods company to enhance its relationship with retailers.

- Establish the timetables for distribution among the outlets.
- Check out the loading patterns.

## VII. CONCLUSION

Data mining tools and techniques have been defined. In particular, we talked about tools and techniques for web data mining. Our warehousing project includes designing for web data mining in order to produce some useful information from the WWW data. We are now looking into the concepts covered in this paper. An overview of data mining, which is the process of extracting useful information from vast amounts of data housed in information repositories, is provided in this paper.

## REFERENCES

[1]. Jiawei Han and Micheline Kamber (2006), Data Mining Concepts and Techniques, published by Morgan Kauffman 2nd

[2]. Dr. Gary Parker, vol 7, 2004, Data Mining: Modules in emerging fields, CD-ROM.

[3]. Crisp-DM 1.0 Step by step Data Mining guide from http://www.crisp-dm.org/CRISPWP-0800.pdf.

[4]. Customer Successes in your industry from http://www.spss.com/success/?source=homepage&hpzone=nav_bar.

[5]. https://www.allbusiness.com/Technology/computer-software-data-management/ 633425-1.html, last retrieved on 15th Aug 2010.

[6]. http://www.kdnuggets.com